# miBLAST and SAGA: Two Scalable NCIBI Bioinformatics Tools

**Jignesh M. Patel**

**Department of EECS**

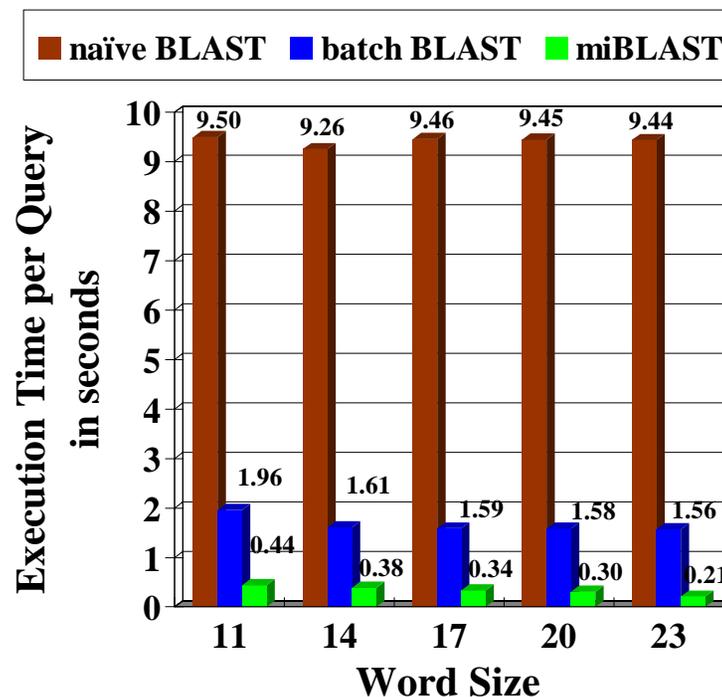**University of Michigan**

**Contact: jignesh@eecs.umich.edu**

**NC BI**

# miBLAST: Scalable BLAST for Batch Workloads

- A common task is to search a large sequence database using a "set" of query sequences.

  – Validation of the Affymetrix probe set against UniGene.

- Running BLAST repeatedly for each query is inefficient.

- Approach: A novel database-inspired "join" algorithm which indexes both the data and the query sets.

- Free download at www.eecs.umich.edu/miblast

- Modifications for MPSS underway (for actual deployment at ISB)

**Query the Affymetrix probe set against Human UniGene**



**miBLAST is 22X faster than BLAST**

National Center for Integrative Biomedical Informatics

2

# SAGA: A Fast and Flexible Graph Matching Tool

- ## Motivation
  - Graph querying is a common requirement for many DBPs.
  - Examples of graph datasets: KEGG, bioNLP, MiMI, …
  - Datasets are noisy/incomplete, so exact matching is inadequate.

- ## Challenge: Graph Matching is a Hard Problem
  - Subgraph isomorphism is NP-complete!
  - *Approximate* Subgraph Matching: Allow approximate matching of node/edge labels, and structural differences.

- ## The database-centric SAGA approach
  - Build an index on small database graph substructures.
  - Use the index to match database and query graph fragments.
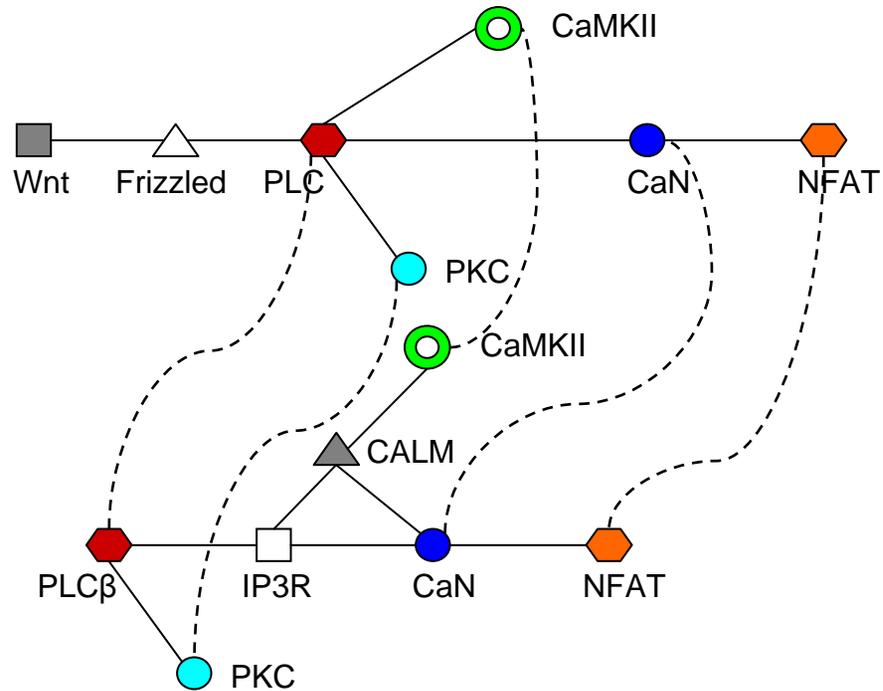  - Assemble larger matches by detecting graph cliques.

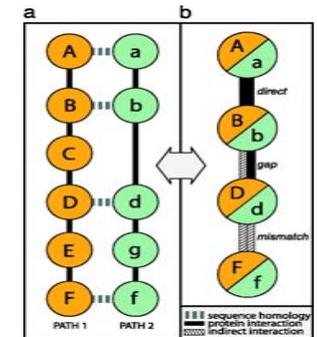# Results: Query KEGG with Wnt/CA2+ Pathway

Query:  Wnt/Ca2+ Signaling

KEGG id: 04310hsa

Match: Calcium Signaling

KEGG id: 04020hsa

CaMKII

Wnt    Frizzled    PLC    CaN    NFAT

PKC

CaMKII

CALM

PLCβ    IP3R    CaN    NFAT

PKC

- • **None of the existing methods can detect this match!**
- • Limitations of Existing Methods:
  - – Gindex & GraphGrep: only perform exact matching
  - – Grafil & PIS: no gap nodes are allowed
  - – PathBlast: only matches paths; edge alignment only tolerates one gap nodes, e.g. (B,D) with (b,d) and (D, F) with (d, f)

**Kelley et. al. PNAS(2003)**

**PathBlast Example**

The Similarity Score: 31.2

Click here to view the

**query:** GR

cd40

cdk6

**Cell Press**

## CDK inhibitor p18(INK4c) is required for the generation of functional plasma cells.

Tourigny MR, Ursini-Siegel J, Lee H, Toellner KM, Cunningham AF, Franklin DS, Ely S, Chen M, Qin XF, Xiong Y, MacLennan IC, Chen-Kiang S.

Department of Pathology, Weill Medical College of Cornell University, 1300 York Avenue, New York, NY 10021, USA.

B cell terminal differentiation is associated with the onset of high-level antibody secretion and cell cycle arrest. Here the cyclin-dependent kinase (CDK) inhibitor p18(INK4c) is shown to be required within B cells for both terminating cell proliferation and differentiation of functional plasma cells. In its absence, B cells hyperproliferate in germinal centers and extrafollicular foci in response to T-dependent antigens but serum antibody titers are severely reduced, despite unimpaired germinal center formation, class switch recombination, variable region-directed hypermutation, and differentiation to antibody-containing plasmacytoid cells. The novel link between cell cycle control and plasma cell differentiation may, at least in part, relate to p18(INK4c) inhibition of CDK6. Cell cycle arrest mediated by p18(INK4C) is therefore requisite for the generation of functional plasma cells.

DR. graphn: GR12668976

FULL TEXT AVAILABLE ONLINE
WILEY InterScience

## p18(INK4c) collaborates with other CDK-inhibitory proteins in the regenerating liver.

Luedde T, Rodriguez ME, Tacke F, Xiong Y, Brenner DA, Trautwein C.

Department of Gastroenterology, Hepatology and Endocrinology, Medizinische Hochschule Hannover, Hannover, Germany.

p18(INK4c) belongs to the family of cyclin-dependent kinase inhibitory proteins that target the cyclin-dependent kinases and inhibit their catalytic activity. The role of p18(INK4c) for cell cycle progression in vivo is characterized poorly. Therefore, we studied the expression and physiologic relevance of p18 in quiescent and proliferating hepatocytes during liver regeneration. For our analysis we used single- (p18[INK4c], p27[KIP1], p21[CIP1/WAF1]), and double-mutant (p18/p21, p18/p27) mice. p18 expression was found in quiescent hepatocytes and a slight up-regulation was evident after partial hepatectomy (PH). p18 knockout animals showed normal cell cycle progression after PH. However, when p18/p21 and p18/p27 double-mutant mice were used, differences in cell cycle progression were evident compared with wild-type (wt) and single knockout animals. In p18/p21 knockout animals, the G1 phase was shortened as evidenced by an earlier onset of cyclin D and proliferating cell nuclear antigen (PCNA) expression and cyclin-dependent kinase (CDK) activation after PH. In contrast, in p18/p27 knockout animals, the G1 phase was unchanged, but the amount of proliferating hepatocytes

# SAGA Demo and Poster